Bias-Variance and Cross Validation

Nipun Batra and teaching staff

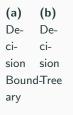
July 18, 2025

IIT Gandhinagar

A Question!

What would be the decision boundary of a decision tree classifier?

Decision Boundary for a tree with depth 1



Decision Boundary for a tree with no depth limit

Are deeper trees always better?

As we saw, deeper trees learn more complex decision boundaries.

Are deeper trees always better?

As we saw, deeper trees learn more complex decision boundaries.

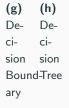
But, sometimes this can lead to poor generalization

An example

Consider the dataset below

Underfitting

Underfitting is also known as *high bias*, since it has a very biased incorrect assumption.



6

Overfitting

Overfitting is also known as *high variance*, since very small changes in data can lead to very different models. Decision tree learned has depth of 10.

Intution for Variance

A small change in data can lead to very different models.

Dataset 1

Dataset 2

Intution for Variance

A Good Fit

As depth increases, train accuracy improves

As depth increases, train accuracy improves As depth increases, test accuracy improves till a point

As depth increases, train accuracy improves
As depth increases, test accuracy improves till a point
At very high depths, test accuracy is not good (overfitting).

Accuracy vs Depth Curve: Underfitting

The highlighted region is the underfitting region. Model is too simple (less depth) to learn from the data.

Accuracy vs Depth Curve: Overfitting

The highlighted region is the overfitting region.

Model is complex (high depth) and hence also learns the anomalies in data.

The highlighted region is the good fit region.

We want to maximize test accuracy while being in this region.

The big question!?

How to find the optimal depth for a decision tree?

The big question!?

How to find the optimal depth for a decision tree?

Use cross-validation!

Our General Training Flow

K-Fold cross-validation: Utilise full dataset for testing

 ${\dots}/{\tt cross-validation-train-test.pdf}$

The Validation Set

Nested Cross Validation

Divide your training set into K equal parts. Cyclically use 1 part as "validation set" and the rest for training. Here K=4

Cross-Validation

Nipun Batra and teaching staff January 12, 2024

IIT Gandhinagar

Nested Cross Validation

Average out the validation accuracy across all the folds Use the model with highest validation accuracy

Next time: Ensemble Learning

- How to combine various models?
- Why to combine multiple models?
- How can we reduce bias?
- How can we reduce variance?